
Optimistic planning for question selection

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In many online educational environments, a common problem is how to select
2 the best questions to give to a user in order to maximize their learning over the
3 limited period of time they interact with the platform. We consider the problem
4 of determining the best policy, or sequence of questions, to give to a user in a
5 fixed period where both the amount of time they spend on each question and the
6 benefit they gain from it are stochastic (but can be simulated). This problem is,
7 in essence, the stochastic knapsack problem, a PSPACE-hard problem that has
8 been studied extensively in operations research. It is desirable that any solution
9 algorithm produces sequences of question that depend on the amount of time that
10 the user spent on the preceding questions. We propose a new adaptive algorithm
11 for the stochastic knapsack problem that uses techniques from multi-armed bandits
12 to explore only potentially good policies. Our algorithm is adaptive to the amount
13 of remaining time and can be shown to obtain a solution within ϵ of the optimal.
14 We then discuss the use of this algorithm in an online education environment.

15 1 Introduction

16 In online education, one key challenge is to provide students with personalized sequences of questions
17 that make the best use of the users time to maximize learning. Imagine an app that a user interacts
18 with for 15 minutes a day in order to learn a new language or skill. The app must decide which
19 questions it can give the student in this time period in order to maximize their learning. However,
20 the amount of time the student will take to complete each question and the benefit they will gain
21 from doing so are both stochastic. This stochasticity makes the problem considerably more difficult.
22 In order to deal with it, we assume we are able to simulate accurately user time and reward from
23 completing each question. Much work has focused on obtaining good predictive models of student
24 performance and question time (see for example, [7, 11, 17]), and we assume we have access to such
25 models at an individual student level (since this leads to better decision making [13]).

26 The stochastic knapsack problem [8], is a classic resource allocation problem that consists of selecting
27 a subset of items to place into a knapsack of a given capacity. Placing each item in the knapsack
28 consumes a random amount of the capacity and provides a stochastic reward. Our problem of
29 selecting which questions to give to a user in a 15 minute exercise can be thought of as the stochastic
30 knapsack problem. Each question (item) will take a random amount of time (size) and improve the
31 student's knowledge in a stochastic manner (reward). To make optimal use of the available time the
32 app needs to track the progress of the user and adjust accordingly. Once an item is placed in the
33 knapsack, we assume we observe its realized size and can use this to make future decisions. This
34 enables us to consider adaptive or closed loop strategies, which will generally perform better [9] than
35 open loop strategies in which the schedule is invariant of the remaining time.

36 For our purposes, it is desirable to have methods for the stochastic knapsack problem that can make
37 use of all available resources and adapt with the remaining capacity. One manner of obtaining such
38 adaptive solutions is to model the problem as a decision tree as discussed in [9]. However, in all

39 but the simplest cases, this decision tree will be too large to search exhaustively. We propose using
40 ideas from optimistic planning [4, 16] to significantly accelerate the tree search by only considering
41 policies with high upper confidence bounds. Most optimistic planning algorithms were developed for
42 discounted mdps and as such rely on discount factors to limit future reward which are not present
43 in the stochastic knapsack problem. Furthermore, in our problem, the random variable representing
44 state transitions also provides us with information on the future rewards. To avoid discount factors
45 and to use the transition information, we work with confidence bounds that incorporate estimates of
46 the remaining capacity and use these estimates to determine how many samples we need. In order
47 to do this, we need techniques that can deal with weak dependencies and that give us confidence
48 regions that hold simultaneously for multiple sample sizes. We therefore combine Doob’s martingale
49 inequality with Azuma-Hoeffding bounds to create high probability bounds. Following the optimistic
50 planning approach, we use these bounds to develop an algorithm that adapts to the complexity of the
51 problem instance: it is guaranteed to find an ϵ -good approximation independent of how difficult the
52 problem is and, if the problem instance is easy to solve, it expands only a moderate sized tree.

53 1.1 Related work

54 Finding exact solutions to the simpler deterministic knapsack problem, in which item weights and
55 rewards are deterministic, is known to be NP-hard and it has been stated that the stochastic knapsack
56 problem is PSPACE-hard [9]. Therefore, most work on the stochastic variant of the problem has
57 focused on approximations. The state-of-the-art approaches to the stochastic knapsack problem
58 where the reward and resource consumption distributions are known, were introduced in [9] where
59 the authors introduced a heuristic that is adaptive and comes within a factor of a 1/3rd of the best
60 total reward. The heuristic groups the available items into small and large items and fills the knapsack
61 exclusively with items of one of the two groups. The strategy for small items is non-adaptive and
62 orders these items according to their reward - consumption ratio, placing items into the knapsack
63 according to this ordering. For the large items, a decision tree is built to some predefined depth d and
64 an exhaustive search for the best policy in that decision tree is performed.

65 Optimistic planning was developed for tree search in large deterministic [12] and stochastic (both
66 open [3] and closed [4] loop) systems. The general idea is to use the upper confidence principle of
67 the UCB algorithm for multi-armed bandits [1] to expand a tree. This is achieved by expanding nodes
68 that have the potential to lead to good solutions through using bounds that take into account both the
69 reward received in getting to a node and the reward that could be obtained after moving on from that
70 node. [16] use optimistic planning in discounted MDPs, requiring only a generative model of the
71 rewards and transitions. Instead of the UCB algorithm, their work relies on the best arm identification
72 algorithm of [10]. Optimistic planning algorithms are used to return a near optimal first action and
73 are then rerun to select the next action. In our case, the decision tree is a good approximation to
74 the entire problem so we can output a near-optimal policy. In our problem, the state transitions
75 (size of item/time taken to answer question) provide information about future rewards and so should
76 be considered when defining the high confidence bounds. Furthermore, our algorithm iteratively
77 builds confidence bounds which are used to determine whether it is necessary to sample more thus
78 making more effective use of resources. One would imagine that the StOP algorithm from [16] could
79 be easily adapted to the stochastic knapsack problem. However, the assumptions required for this
80 algorithm to terminate are too strong for it to be considered a feasible algorithm for our problem.

81 In the education literature, there has been some work done on the question selection problem. [6]
82 propose to use a non-stationary multi-armed bandit algorithm to select questions for students, without
83 the time constraints that are present in our problem. [15] consider the use of POMDPs to determine
84 which type of action to take next based on previous success but focus on the choice between types of
85 actions (such as videos, quizzes or questions with feedback) rather than the specific question itself. [5]
86 also use reinforcement learning algorithms to work towards a similar aim.

87 1.2 Our contribution

88 Our main contributions are the confidence bounds in Lemma 1 and Proposition 2 that allow us to
89 simultaneously estimate remaining capacity and reward, with guarantees that hold uniformly over
90 multiple sample sizes; Proposition 3, which shows that we can avoid discount based arguments
91 and still return an adaptive policy with value within ϵ of the optimal policy with high probability

92 while using varying capacity estimates; and, primarily, our algorithm OpStoK whose use to provide
 93 instructional policies for education software will be discussed in Section 5.

94 2 Problem formulation

95 In this section, we formally define our problem. Note that since our problem is essentially the
 96 stochastic knapsack problem, we define it in terms of the knapsack definitions and describe it using
 97 items, sizes and budgets, rather than questions, time taken and time limits.

98 We consider the problem of selecting a subset of items from a set of K items, I , to place into a
 99 knapsack of capacity B where each item can be played at most once. For each item $i \in I$, let C_i and
 100 R_i be bounded random variables defined on a joint probability space (Ω, \mathcal{A}, P) which represent the
 101 size and reward of item i . It is assumed that we can simulate from the generative model of (R_i, C_i)
 102 for all $i \in I$ and we will use lower case c_i and r_i , to denote realizations of the random variables.
 103 We assume that the random variables (R_i, C_i) are independent of (R_j, C_j) for all $i, j \in I, i \neq j$.
 104 Further, it is believed that item sizes and rewards do not change dependent on the other items in the
 105 knapsack. We assume the problem is non-trivial, in the sense that it is not possible to fit all items in
 106 the knapsack at once. If we place an item i in the knapsack and the consumption C_i is strictly greater
 107 than the remaining capacity then we gain no reward for this item. Our final important assumption is
 108 that there exists some non-decreasing function $\Psi(\cdot)$, satisfying $\lim_{b \rightarrow 0} \Psi(b) = 0$ and $\Psi(B) < \infty$,
 109 such that the reward that can be achieved with budget b is upper bounded by $\Psi(b)$.

110 Representing the stochastic knapsack problem as a tree requires that all item sizes take discrete values.
 111 While in this work, it will generally be assumed that this is the case, in some problem instances,
 112 continuous item sizes need to be discretized. In this case, let v^* be the optimal value of the best
 113 policy and let ξ^* be the corresponding discretization error. Then $\Psi(\xi^*)$ is an upper bound on the
 114 extra reward that could be gained from the space lost due to discretization. For discrete sizes, we
 115 assume there are s possible values the random variable can take and that there exists a value $\theta > 0$
 116 such that $C_i \geq \theta$ for all $i \in I$.

117 2.1 Planning trees and policies

118 The stochastic knapsack problem can be thought of as a planning tree with the initial empty state as
 119 the root at level 0. Each node on an even level is an *action* node and its children represent placing an
 120 item in the knapsack. The nodes on odd levels are *transition* nodes with children representing item
 121 sizes. We define a *policy* Π as a finite subtree where each action node has at most one child and each
 122 transition node has s children. The *depth* of a policy Π , $d(\Pi)$, is defined as the number of transition
 123 nodes in any realization of the policy (where each transition node has one child), or equivalently,
 124 the number of items. Let $d^* = \lfloor B/\theta \rfloor$ be the maximal depth of any policy. For any $1 \leq d \leq d^*$, the
 125 number of policies of depth d is

$$N_d = \prod_{i=0}^{d-1} (K - i)^{s^i} \quad (1)$$

126 where $K = |I|$ is the number of items, and s the number of discrete sizes.

127 We define a *child* policy, Π' , of a policy Π as a policy that follows Π up to depth $d(\Pi)$ then plays
 128 additional items and has depth $d(\Pi') = d(\Pi) + 1$. In this setting, Π is the *parent* policy of Π' . A
 129 policy is said to be *incomplete* if the remaining budget allows for another item to be inserted into the
 130 knapsack (see Section 4 for a formal definition).

131 The *value* of a policy Π can be defined as the cumulative expected reward obtained by playing items
 132 according to Π , $V_\Pi = \sum_{t=1}^T E[R_{i_t}]$ where i_t is the t -th item chosen by Π . Let \mathcal{P} be the set of all
 133 policies, then define the *optimal policy* as $\Pi^* = \arg \max_{\Pi \in \mathcal{P}} V_\Pi$, and corresponding *optimal value*
 134 as $v^* = \max_{\Pi \in \mathcal{P}} V_\Pi$. Our algorithm returns an ϵ -optimal policy with value $v^* - \epsilon$. For any policy
 135 Π , we define a *sample* of Π as follows. The first item of any policy is fixed so we take a sample of
 136 the reward and size from the generative model of that item. We then use Π to tell us which item to
 137 sample next (based on the size of the previous item) and sample the reward and size of that item. This
 138 continues until the policy finishes or the cumulative size of the selected items exceeds B .

139 3 High confidence bounds

140 In this section, we develop confidence bounds for the value of a policy. Observe that a policy Π
 141 need not consume all available budget, in fact our algorithm will construct iteratively longer policies
 142 starting from the shortest policies of playing a single item. Consequently, we are also interested in
 143 R_{Π}^+ , the expected maximal reward that can be obtained after playing according to policy Π until
 144 all the budget is consumed. Let B_{Π} be a random variable representing the remaining budget after
 145 playing according to a policy Π . Our assumptions guarantee that there exists a function Ψ such that
 146 $R_{\Pi}^+ \leq E\Psi(B_{\Pi})$. We define V_{Π}^+ to be the maximal expected value of any continuation of policy Π so
 147 $V_{\Pi}^+ = V_{\Pi} + R_{\Pi}^+ \leq V_{\Pi} + E\Psi(B_{\Pi})$.

148 From m samples of the reward of policy Π , we estimate the value of Π as $\overline{V}_{\Pi m} =$
 149 $\frac{1}{m} \sum_{j=1}^m \sum_{d=1}^{d(\Pi)} r_{i(d)}^{(j)}$, where $r_{i(d)}^{(j)}$ is the reward of item $i(d)$ chosen at depth d of sample j .
 150 However, our real interest is in the value of V_{Π}^+ since we wish to identify the policy with
 151 greatest reward when continued until the budget is exhausted. From Hoeffding's inequality,

$$152 P\left(\left|\overline{V}_{\Pi m_1} - V_{\Pi}^+\right| > E\Psi(B_{\Pi}) + \sqrt{\frac{\Psi(B)^2 \log(2/\delta)}{2m}}\right) \leq \delta. \text{ This bound depends on the quantity}$$

153 $E\Psi(B_{\Pi})$ which is typically not known. Furthermore, our algorithm will work by sampling $\Psi(B_{\Pi})$
 154 until we are confident enough that it is small or large. As such, it introduces weak dependencies
 155 into the sampling process and we need confidence bounds that will hold simultaneously for multiple
 156 numbers of samples of the remaining capacity, m_2 . Hence, we work with martingale techniques
 157 and use Azuma-Hoeffding like bounds [2], similar to the technique used in [14]. Specifically, in
 158 Lemma 3 (supplementary material), we use Doob's maximal inequality and a peeling argument to get
 159 Azuma like bounds for the maximal deviation of the sample mean from the expected value under
 160 boundedness assumptions. Assuming we sample the reward m_1 times and remaining capacity of a
 161 policy $m_2 \leq n$ times, the following key result holds.

162 **Proposition 1** *The Algorithm BoundValueShare (Algorithm 2) returns confidence bounds,*

$$L(V_{\Pi}^+) = \overline{V}_{\Pi m_1} - c_1$$

$$U(V_{\Pi}^+) = \overline{V}_{\Pi m_1} + \overline{\Psi(B_{\Pi})}_{m_2} + c_1 + c_2$$

163 *which hold with probability $1 - \delta_1 - \delta_2$, where $c_1 = \sqrt{\frac{\Psi(B)^2 \log(2/\delta_1)}{2m_1}}$, $c_2 = 2\Psi(B) \sqrt{\frac{1}{m_2} \log\left(\frac{8n}{\delta m_2}\right)}$.*

164 This upper bound depends on n , the maximum number of samples of $\Psi(B_{\Pi})$. For any policy Π , the
 165 minimum width of a confidence interval of $\Psi(B_{\Pi})$ required by the algorithm BoundValueShare
 166 when run with precision parameter ϵ is $\epsilon/4$. Hence, taking,

$$n = \left\lceil \frac{16^2 \Psi(B)^2 \log(8/\delta)}{\epsilon^2} \right\rceil, \quad (2)$$

167 ensures that for all policies, $2c_2 \leq \epsilon/4$ when $m_2 = n$. As discussed in Section 4, this is a necessary
 168 condition for the termination of the algorithm.

169 4 Algorithm

170 The process of sampling the reward of a policy involves sampling item sizes to decide which item
 171 to play next. We propose to make better use of all available data by using the samples of item sizes
 172 to calculate $U(\Psi(B_{\Pi}))$. Our algorithm, OpStoK, will then use the tight upper bound $U(\Psi(B_{\Pi}))$ in
 173 the bound future rewards, $U(V_{\Pi}^+)$. We also pool samples of the reward and size of items across
 174 policies, thus reducing the number of calls to the generative model. OpStoK also benefits from an
 175 alternative sampling method that reduces sample complexity and ensures that an entire ϵ -optimal
 176 policy is returned when the algorithm stops (line 5, Algorithm 1). This is achieved by using the bound
 177 in Proposition 1 and n as defined in (2).

178 In the main algorithm OpStoK (Algorithm 1), is very similar to StOP-K from [16] with the
 179 key differences appearing in the sampling and construction of confidence bounds, defined in
 180 BoundValueShare, that ensure the algorithm converges. OpStoK proceed by maintaining a set

Algorithm 1: OpStoK ($I, \delta_{0,1}, \delta_{0,2}, \epsilon$)

Initialization: ACTIVE = \emptyset

- 1 **forall the** $i \in I$ **do**
- 2 Π_i = policy consisting of just playing item i .
- 3 $d(\Pi_i) = 1$
- 4 $\delta_{1,1} = \frac{\delta_{0,1}}{d^*} N_1^{-1}$ $\delta_{1,2} = \frac{\delta_{0,2}}{d^*} N_1^{-1}$
- 5 $(L(V_{\Pi_i}^+), U(V_{\Pi_i}^+)) = \text{BoundValueShare}(\Pi_i, \delta_{0,1}, \delta_{0,2}, \mathcal{S}^*, \epsilon)$
- 6 ACTIVE = ACTIVE $\cup \{\Pi_i\}$.
- 7 **end**
- 8 **for** $t = 1, 2, \dots$ **do**
- 9 $\Pi_t^{(1)} = \arg \max_{\Pi \in \text{ACTIVE}} U(V_{\Pi}^+)$
- 10 $\Pi_t^{(2)} = \arg \max_{\Pi \in \text{ACTIVE} \setminus \{\Pi_t^{(1)}\}} U(V_{\Pi}^+)$
- 11 **if** $L(V_{\Pi_t^{(1)}}^+) + \epsilon \geq \max_{\Pi \in \text{ACTIVE}} U(V_{\Pi}^+)$ **then**
- 12 **Stop:** $\Pi^* = \Pi_t^{(1)}$;
- 13 $\Pi_t = \Pi_t^{(a^*)}$, where $a^* = \arg \max_{a \in \{1,2\}} U(\Psi(B_{\Pi_t^{(a)}}))$
- 14 ACTIVE = ACTIVE $\setminus \{\Pi_t\}$
- 15 **forall the children** Π' **of** Π_t **do**
- 16 $d(\Pi') = d(\Pi_t) + 1$
- 17 $\delta_1 = \frac{\delta_{0,1}}{d^*} N_{d(\Pi')}^{-1}$ and $\delta_2 = \frac{\delta_{0,2}}{d^*} N_{d(\Pi')}^{-1}$
- 18 $(L(V_{\Pi'}^+), U(V_{\Pi'}^+)) = \text{BoundValueShare}(\Pi', \delta_1, \delta_2, \mathcal{S}^*, \epsilon)$
- 19 ACTIVE = ACTIVE $\cup \{\Pi'\}$
- 20 **end**
- 21 **end**

181 of ‘active’ policies. As in [16] and [10], at each time step t , a policy, Π_t to expand is chosen by
182 comparing the upper confidence bounds of the two best active policies. We select the policy with most
183 uncertainty in the bounds since we want to be confident enough in our estimates of the near-optimal
184 policies to say that the policy we ultimately select is better (see Figure 1). Once we have selected
185 a policy, Π_t , if the stopping criteria is not met, we replace Π_t in the set of active policies with all
186 its children. For each child policy, we use BoundValueShare to bound its reward. In order for all
187 our bounds to hold simultaneously with probability greater than $1 - \delta_{0,1} - \delta_{0,2}$, BoundValueShare
188 must be called with parameters

$$\delta_{d,1} = \frac{\delta_{0,1}}{d^*} N_{d(\Pi)}^{-1} \quad \text{and} \quad \delta_{d,2} = \frac{\delta_{0,2}}{d^*} N_{d(\Pi)}^{-1} \quad (3)$$

189 where N_d is the number of policies of depth d as given in (1). Our algorithm, OpStoK is given in
190 Algorithm 1. The algorithm relies on BoundValueShare and subroutines, EstimateValue and
191 SampleBudget, which sample the reward and budget of policies.

192 In BoundValueShare, we use samples of both item size and reward to bound the value of a policy.
193 We define upper and lower bounds on the value of any extension of a policy Π as,

$$\begin{aligned} U(V_{\Pi}^+) &= \overline{V_{\Pi} m_1} + \overline{\Psi(B_{\Pi})}_{m_2} + c_1 + c_2, \\ L(V_{\Pi}^+) &= \overline{V_{\Pi} m_1} - c_1, \end{aligned}$$

194 with c_1 and c_2 as in Proposition 1. It is also possible to define upper and lower bounds on $\Psi(B_{\Pi})$ with
195 m_2 samples and confidence δ_2 . From this, we can formally define a *complete* policy as a policy Π with
196 $U(B_{\Pi}) = \overline{\Psi(B_{\Pi})}_{m_2} + c_2 \leq \frac{\epsilon}{2}$. For complete policies, since there is very little capacity left, it is more
197 important to get tight confidence bounds on the value of the policy. Hence, in BoundValueShare,
198 we sample the remaining budget policy as much as is necessary to conclude whether the policy is
199 complete or not. As soon as we realize we have a complete policy ($U(B_{\Pi}) \leq \epsilon/2$), we sample the
200 value sufficiently to get a confidence interval of width less than ϵ . Then, when it comes to choosing
201 an optimal policy to return, the confidence intervals of all complete policies will be narrow enough
202 for this to happen. This is appropriate since, pre-specifying the number of samples may not lead

Algorithm 2: BoundValueShare($\Pi, \delta_1, \delta_2, S^*, \epsilon$)

Input: Π : policy; δ_1 : probability capacity confidence bound fails; δ_2 : probability reward confidence bound fails; S^* : observed samples for all items; ϵ : tolerated approximation error.

Initialization: For all $i \in I$, let $\mathcal{S}_i = S_i^*$

```

1 Set  $m_2 = 1$  and  $(\psi_1, \mathcal{S}) = \text{SampleBudget}(\Pi, \mathcal{S})$ 
  /* draw a sample of the remaining budget */
2  $\overline{\Psi}(B_\Pi)_{m_2} = \frac{1}{m_2} \sum_{j=1}^{m_2} \psi_j$ 
3  $U(\Psi(B_\Pi)) = \overline{\Psi}(B_\Pi)_{m_2} + 2\Psi(B) \sqrt{\frac{1}{m_2} \log\left(\frac{8n}{\delta m_2}\right)}$ ,
   $L(\Psi(B_\Pi)) = \overline{\Psi}(B_\Pi)_{m_2} - 2\Psi(B) \sqrt{\frac{1}{m_2} \log\left(\frac{8n}{\delta m_2}\right)}$ 
  /* calculate upper and lower bounds on the remaining budget */
4 if  $U(\Psi(B_\Pi)) \leq \frac{\epsilon}{2}$  then  $m_1 = \left\lceil \frac{8\Psi(B)^2 \log(2/\delta_1)}{\epsilon^2} \right\rceil$ ;
5 else if  $L(\Psi(B_\Pi)) \geq \frac{\epsilon}{4}$  then
6    $m_1 = \left\lceil \frac{1}{2} \frac{\Psi(B)^2 \log(2/\delta_1)}{u(\Psi(B))^2} \right\rceil$ 
7 else
8   Set  $m_2 = m_2 + 1$ ,  $(\psi_{m_2}, \mathcal{S}) = \text{SampleBudget}(\Pi, \mathcal{S})$  and go back to 2
9  $\overline{V}_{\Pi m_1} = \text{EstimateValue}(\Pi, m_1)$ 
10  $L(V_\Pi^+) = \overline{V}_{\Pi m_1} - \sqrt{\frac{\Psi(B)^2 \log(2/\delta_1)}{2m_1}}$ 
11  $U(V_\Pi^+) = \overline{V}_{\Pi m_1} + \overline{\Psi}(B_\Pi)_{m_2} + \sqrt{\frac{\Psi(B)^2 \log(2/\delta_1)}{2m_1}} + 2\Psi(B) \sqrt{\frac{1}{m_2} \log\left(\frac{8n}{\delta m_2}\right)}$ 
12 return  $(L(V_\Pi^+), U(V_\Pi^+))$ 

```

203 to confidence bounds tight enough to select an ϵ -optimal policy. If a complete policy is chosen as
 204 $\Pi_t^{(1)}$ in OpStoK, for some t , the algorithm will stop and this policy will be returned. For this to
 205 happen, we also need the stopping criterion to be checked before selecting a policy to expand. Note
 206 that in BoundValueShare, the reward and remaining budget must be sampled separately as we are
 207 considering closed-loop planning so the item chosen may depend on the size of the previous item,
 208 and hence the reward will depend on the instantiated item sizes. In line 6 of BoundValueShare, the
 209 number of samples of the reward, m_1 , is defined to ensure that the uncertainty in the estimate of
 210 V_Π is less than $u(\Psi(B)) = \min\{U(\Psi(B_\Pi)), \Psi(B)\}$, since a natural upper bound for the reward is
 211 $\Psi(B)$. In the other case (when $U(\Psi(B_\Pi)) \leq \epsilon/2$), we define it to ensure the confidence bounds are
 212 tight enough.

213 OpStoK, considerably reduces the number of calls to the generative model by creating sets \mathcal{S}_i^* of
 214 samples of the reward and size of each item $i \in I$. When it is necessary to sample the reward and size
 215 of an item for the evaluation of a policy, we sample without replacement from \mathcal{S}_i^* , until $|\mathcal{S}_i^*|$ samples
 216 have been taken. At this point new calls to the generative model are made and the new samples added
 217 to the sets for use by future policies. We denote by S^* the collection of all sets \mathcal{S}_i^* .

218 4.1 Analysis

219 We state the following result guaranteeing the performance of OpStoK

220 **Proposition 2** *With probability at least $(1 - \delta_{0,1} - \delta_{0,2})$, the algorithm OpStoK returns an action*
 221 *with value at least $v^* - \epsilon$ for $\epsilon > 0$.*

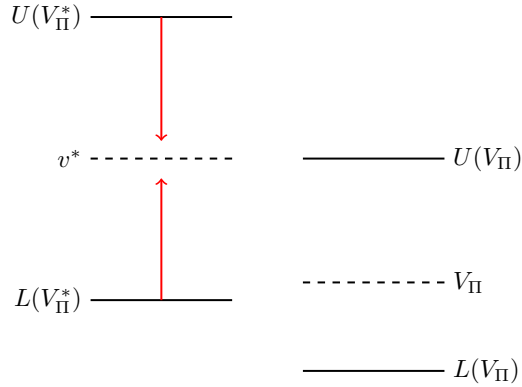


Figure 1: Example of where just looking at the optimistic policy might fail: If we always play the optimistic policy, as $U(V_{\Pi^*}^+) \geq U(V_{\Pi}^+)$, we will always play Π^* so the confidence bounds on Π will not shrink. This means that $L(V_{\Pi^*}^+)$ will never be (epsilon) greater than the best alternative upper bound so there will not be enough confidence to conclude we have found the best policy.

222 5 Applications to online education

223 Our algorithm as discussed here has been in a mainly theoretical framework. However, we believe that
 224 it can have serious practical impact in the field of online education. For the problem of determining
 225 which question to give to a user in a fixed time frame, the tree based structure to question selection
 226 feels very natural. It allows for pre-requisite exercises to be included very easily, which along with
 227 producing more pedagogically sound policies will also improve computational efficiency as the
 228 search space will be reduced at the policy expansion stage. With this in mind, it may also be possible
 229 to incorporate ideas from the Zone of Proximal Development into the algorithm. Additionally, it
 230 would be interesting to consider reward which change depending on which questions have previously
 231 been asked. OpStoK provides an intuitive manner of dealing with the large decision trees that can
 232 arise when considering the problem of question selection in online education. Our algorithm will
 233 only evaluate potentially optimal policies and can be run offline (provided a good model of student
 234 performance and time are available) to produce educational policies that adapt to the remaining time
 235 and provide a near optimal learning experience for the limited period of time the user has to spend on
 236 the app or platform.

237 6 Conclusion

238 In this paper we have presented a new algorithm OpStoK that applies the optimistic planning strategy
 239 to the stochastic knapsack problem. This algorithm is directly motivated by problem of selecting
 240 an adaptive sequence of questions to give to a student in an education app. At present, we have
 241 provided largely theoretical results so in future it would be good to investigate the performance
 242 of our algorithm in a real education environment. While it is suspected that the algorithm will be
 243 computationally efficient as it reduces the number of policies that have to be explored, it would be
 244 interesting to run some experiments both to investigate its practical performance, both in artificial
 245 environments and the true educational setting for which it was designed.

246 **References**

- 247 [1] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem.
248 *Machine Learning*, 47(2-3):235–256, 2002.
- 249 [2] K. Azuma. Weighted sums of certain dependent random variables. *Tohoku Math. J. (2)*, 1967.
- 250 [3] S. Bubeck and R. Munos. Open loop optimistic planning. In *Conference on Learning Theory*,
251 pages 477–489, 2010.
- 252 [4] L. Busoniu and R. Munos. Optimistic planning for markov decision processes. In *15th*
253 *International Conference on Artificial Intelligence and Statistics*, volume 22, pages 182–189,
254 2012.
- 255 [5] M. Chi, K. VanLehn, D. Litman, and P. Jordan. Empirically evaluating the application of
256 reinforcement learning to the induction of effective and adaptive pedagogical strategies. *User*
257 *Modeling and User-Adapted Interaction*, 21(1-2):137–180, 2011.
- 258 [6] B. Clement, D. Roy, P.-Y. Oudeyer, and M. Lopes. Multi-armed bandits for intelligent tutoring
259 systems. *arXiv preprint arXiv:1310.3174*, 2013.
- 260 [7] A. T. Corbett and J. R. Anderson. Knowledge tracing: Modeling the acquisition of procedural
261 knowledge. *User modeling and user-adapted interaction*, 4(4):253–278, 1994.
- 262 [8] G. B. Dantzig. Discrete-variable extremum problems. *Operations Research*, 5(2):266–288,
263 1957.
- 264 [9] B. C. Dean, M. X. Goemans, and J. Vondrák. Approximating the stochastic knapsack problem:
265 The benefit of adaptivity. *Mathematics of Operations Research*, 33(4):945–964, 2008.
- 266 [10] V. Gabillon, M. Ghavamzadeh, and A. Lazaric. Best arm identification: A unified approach
267 to fixed budget and fixed confidence. In *Advances in Neural Information Processing Systems*,
268 pages 3212–3220, 2012.
- 269 [11] R. K. Hambleton. *Fundamentals of item response theory*. Sage publications, 1991.
- 270 [12] J.-F. Hren and R. Munos. Optimistic planning of deterministic systems. In *Recent Advances in*
271 *Reinforcement Learning*, pages 151–164. Springer, 2008.
- 272 [13] J. I. Lee and E. Brunskill. The impact on individualizing student models on necessary practice
273 opportunities. *International Educational Data Mining Society*, 2012.
- 274 [14] V. Perchet, P. Rigollet, S. Chassang, and E. Snowberg. Batched bandit problems. *The Annals of*
275 *Statistics*, 44(2):660–681, 2016.
- 276 [15] A. N. Rafferty, E. Brunskill, T. L. Griffiths, and P. Shafto. Faster teaching via pomdp planning.
277 *Cognitive Science*, 2015.
- 278 [16] B. Szörényi, G. Kedenburg, and R. Munos. Optimistic planning in markov decision processes
279 using a generative model. In *Advances in Neural Information Processing Systems*, pages
280 1035–1043, 2014.
- 281 [17] W. J. van der Linden. A lognormal model for response times on test items. *Journal of*
282 *Educational and Behavioral Statistics*, 31(2):181–204, 2006.
- 283 [18] D. Williams. *Probability with martingales*. Cambridge university press, 1991.

284 **Proof of Proposition 1**

285 The proof of Proposition 1 relies on the following key result.

286 **Lemma 3** Let $\{Z_m\}_{m=1}^\infty$ be a martingale with Z_m defined on the filtration \mathcal{F}_m , $E[Z_m] = 0$ and
 287 $|Z_m - Z_{m-1}| \leq d$ for all m where $Z_0 = 0$. Then,

$$P\left(\exists m \leq n; \frac{Z_m}{m} \geq 2d^2 \sqrt{\frac{2}{m} \log\left(\frac{n}{m} \frac{4}{\delta}\right)}\right) \leq \delta$$

288 *Proof:* The proof is similar to that of Lemma B.1 in [14] and will make use of the following standard
 289 results:

290 **Theorem 4** Doob's maximal inequality: let Z be a non-negative submartingale. Then for $c > 0$,

$$P\left(\sup_{k \leq n} Z_k \geq c\right) \leq \frac{E[Z_n]}{c}.$$

291 *Proof:* See, for example, [18], Theorem 14.6, page 137. □

292 **Lemma 5** Let Z_n be a martingale such that $|Z_i - Z_{i-1}| \leq d_i$ for all i with probability 1. Then, for
 293 $\lambda > 0$,

$$E[e^{\lambda Z_n}] \leq e^{\frac{\lambda^2 D^2}{2}},$$

294 where $D^2 = \sum_{i=1}^n d_i^2$.

295 *Proof:* See the proof of the Azuma-Hoeffding inequality in [2]. □

296 Then, for the proof of Lemma 3, we first notice that since $\{Z_m\}_{m=1}^\infty$ is a martingale, by Jensen's
 297 inequality for conditional expectations, it follows that for any $\lambda > 0$,

$$E[e^{\lambda Z_m} | \mathcal{F}_{m-1}] \geq e^{\lambda E[Z_m | \mathcal{F}_{m-1}]} = e^{\lambda Z_{m-1}}.$$

298 Hence, for any $\lambda > 0$, $\{e^{\lambda Z_m}\}_{m=1}^\infty$ is a positive sub-martingale so we can apply Doob's maximal
 299 inequality (Theorem 4) to get

$$P\left(\sup_{m \leq n} Z_m \geq c\right) = P\left(\sup_{m \leq n} e^{\lambda Z_m} \geq e^{\lambda c}\right) \leq \frac{E[e^{\lambda Z_n}]}{e^{\lambda c}}.$$

300 Then, by Lemma 5, since $|Z_i - Z_{i-1}| \leq d$ for all i , it follows that

$$P\left(\sup_{m \leq n} Z_m \geq c\right) \leq \frac{E[e^{\lambda Z_n}]}{e^{\lambda c}} \leq \frac{e^{\lambda^2 D^2/2}}{e^{\lambda c}} = \exp\left\{\frac{\lambda^2 D^2}{2} - \lambda c\right\}. \quad (4)$$

301 Minimizing the right hand side with respect to λ gives $\hat{\lambda} = \frac{c}{D^2}$ and substituting this back into (4) we
 302 get,

$$P\left(\sup_{m \leq n} Z_m \geq c\right) \leq \exp\left\{-\frac{c^2}{2D^2}\right\}.$$

303 Then, since we are considering the case where $d_i = d$ for all i , $D^2 = nd^2$ and so,

$$P\left(\sup_{m \leq n} Z_m \geq c\right) \leq \exp\left\{-\frac{c^2}{2nd^2}\right\}.$$

304 Further, if we are interested in $P(\sup_{k \leq m \leq n} Z_m \geq c)$, we can redefine the indices's to get

$$P\left(\sup_{k \leq m \leq n} Z_m \geq c\right) = P\left(\sup_{m' \leq n-k+1} Z_m \geq c\right) \leq \exp\left\{-\frac{c^2}{2(n-k+1)d^2}\right\}. \quad (5)$$

305 We then define $\varepsilon_m = 2d\sqrt{\frac{1}{m} \log\left(\frac{n}{m} \frac{8}{\delta}\right)}$ and use a peeling argument similar to that in Lemma B.1 of
 306 [14] to get

$$\begin{aligned}
 P\left(\exists m \leq n; \frac{Z_m}{m} \geq \varepsilon_m\right) &\leq \sum_{t=0}^{\lfloor \log_2(n) \rfloor + 1} P\left(\bigcup_{m=2^t}^{2^{t+1}-1} \left\{\frac{Z_m}{m} \geq \varepsilon_m\right\}\right) && \text{(by union bound)} \\
 &\leq \sum_{t=0}^{\lfloor \log_2(n) \rfloor + 1} P\left(\bigcup_{m=2^t}^{2^{t+1}-1} \left\{\frac{Z_m}{m} \geq \varepsilon_{2^{t+1}}\right\}\right) && \text{(since } \varepsilon_m \text{ decreasing in } m) \\
 &\leq \sum_{t=0}^{\lfloor \log_2(n) \rfloor + 1} P\left(\bigcup_{m=2^t}^{2^{t+1}-1} \{Z_m \geq 2^t \varepsilon_{2^{t+1}}\}\right) && \text{(as } m \geq 2^t) \\
 &\leq \sum_{t=0}^{\lfloor \log_2(n) \rfloor + 1} \exp\left\{-\frac{(2^t \varepsilon_{2^{t+1}})^2}{2^{t+1} d^2}\right\} && \text{(from (5))} \\
 &\leq \sum_{t=0}^{\lfloor \log_2(n) \rfloor + 1} \frac{2^{t+1} \delta}{8n} && \text{(substituting } \varepsilon_{2^{t+1}}) \\
 &\leq \frac{2^{\log_2(n)+3} \delta}{8n} = \delta. && \text{(since } \sum_{i=1}^k 2^i = 2^{k+1} - 1)
 \end{aligned}$$

307

□

308 We are now able to prove the following proposition.

309 **Proposition 6** (Proposition 1 in main text) *The Algorithm BoundValueShare (Algorithm 2) returns*
 310 *confidence bounds,*

$$\begin{aligned}
 L(V_\Pi) &= \overline{V_{\Pi m_1}} - \sqrt{\frac{\Psi(B)^2 \log(2/\delta_1)}{2m_1}} \\
 U(V_\Pi) &= \overline{V_{\Pi m_1}} + \overline{\Psi(B_\Pi)_{m_2}} + \sqrt{\frac{\Psi(B)^2 \log(2/\delta_1)}{2m_1}} + 2\Psi(B) \sqrt{\frac{1}{m} \log\left(\frac{8n}{\delta m}\right)}
 \end{aligned}$$

311 which hold with probability $1 - \delta_1 - \delta_2$.

312 *Proof:*

313 We then begin by noting that our samples of item size are dependent because we evaluate in each
 314 iteration a bound based on past samples and we use this bound to decide if we need to continue
 315 sampling or if we can stop. To model this dependence let us introduce a stopping time τ such that
 316 $\tau(\omega) = n$ if our algorithm exits the loop at n . Consider the sequence

$$\overline{\Psi(B_\Pi)_{1 \wedge \tau}}, \overline{\Psi(B_\Pi)_{2 \wedge \tau}}, \dots$$

317 and define for $m \geq 1$

$$M_m = (m \wedge \tau) (\overline{\Psi(B_\Pi)_{m \wedge \tau}} - E[\Psi(B_\Pi)]) \quad \text{with} \quad M_0 = 0.$$

318 Furthermore, define the filtration $\mathcal{F}_m = \sigma(B_{\Pi,1}, \dots, B_{\Pi,m})$ then for $m \geq 1$

$$E[M_m | \mathcal{F}_{m-1}] = E[M_m | \mathcal{F}_{m-1}, \tau \leq m-1] + E[M_m | \mathcal{F}_{m-1}, \tau > m-1].$$

319 Now

$$E[M_m | \mathcal{F}_{m-1}, \tau \leq m-1] = E[M_{m-1} | \tau \leq m-1].$$

320 and due to independence of the samples $B_{\Pi,1}, \dots, B_{\Pi,m}$

$$\begin{aligned}
& E[M_m | \mathcal{F}_{m-1}, \tau > m-1] \\
&= E[m(\overline{\Psi(B_{\Pi})}_m) - E[\Psi(B_{\Pi})] | \mathcal{F}_{m-1}, \tau > m-1] \\
&= E \left[\sum_{j=1}^{m-1} \Psi(B_{\Pi,j}) + \Psi(B_{\Pi,m}) - mE[\Psi(B_{\Pi})] \middle| \mathcal{F}_{m-1}, \tau > m-1 \right] \\
&= (m-1)E[\overline{\Psi(B_{\Pi})}_{m-1} - E[\Psi(B_{\Pi})] | \mathcal{F}_{m-1}, \tau > m-1] \\
&\quad + E[\Psi(B_{\Pi,m}) - E[\Psi(B_{\Pi})] | \mathcal{F}_{m-1}, \tau > m-1] \\
&= E[M_{m-1} | \tau > m-1] + E[\Psi(B_{\Pi,m})] - E[\Psi(B_{\Pi})] = E[M_{m-1} | \tau > m-1].
\end{aligned}$$

321 Hence, $E[M_m | \mathcal{F}_{m-1}] = M_{m-1}$ and M_m is a martingale with increments $|M_m - M_{m-1}| \leq$
322 $|\Psi(B_{\Pi,m}) - E[\Psi(B_{\Pi})]| \leq \Psi(B)$. We could apply the Azuma-Hoeffding inequality to gain guaran-
323 tees for individual m -values. Alternatively, we can use Lemma 3 to get,

$$P \left(\sup_{m \leq n} \frac{M_m}{m} \geq 2\Psi(B) \sqrt{\frac{1}{m} \log \left(\frac{8n}{\delta m} \right)} \right) \leq \delta_2.$$

324 Using this in conjunction with Azuma bounds on the reward of a policy gives

$$\overline{V_{\Pi m_1}} - c_1 \leq V_{\Pi}^+ \leq \overline{V_{\Pi m_1}} + \overline{\Psi(B_{\Pi})}_{m_2} + c_1 + c_2,$$

325 where $c_1 := \sqrt{\frac{\Psi(B)^2 \log(2/\delta_1)}{2m_1}}$ and $c_2 := 2\Psi(B) \sqrt{\frac{1}{m} \log \left(\frac{8n}{\delta m} \right)}$ and these bounds hold with probabil-
326 ity $1 - \delta_1 - \delta_2$.

327

□